

MobileBERT On RISC-V: Leveraging IREE Compiler and ACE-RVV Extension for Softmax Acceleration



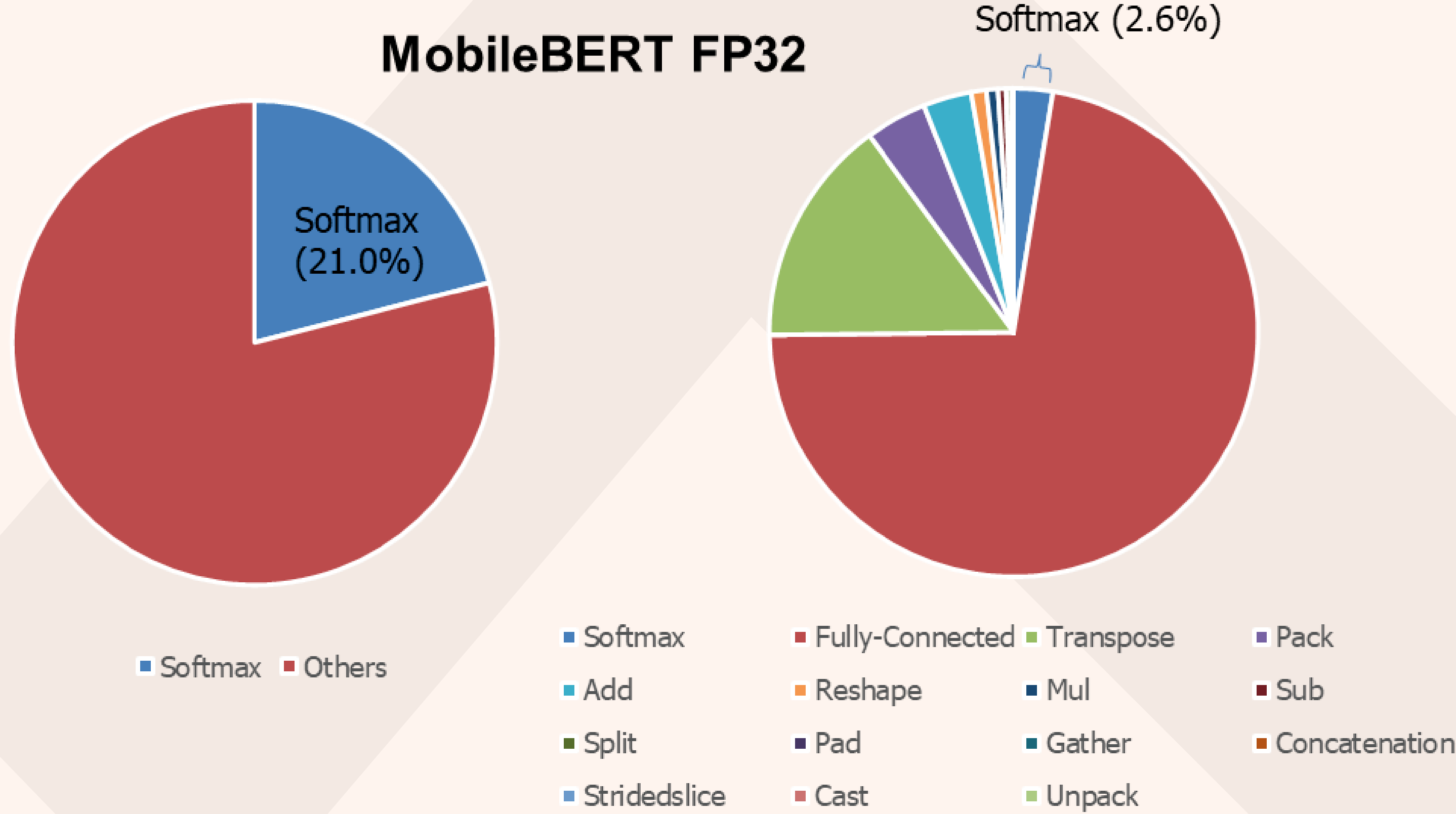
Motivation

► IREE softmax RVV codegen is not as efficient as manually optimized library

► Nonlinear functions based on the exponential function

IREE RVV Codegen

Manually Optimized Library



• Softmax

$$\sigma(\vec{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

• Tanh

$$\tanh x = \frac{\sinh x}{\cosh x} = \frac{e^x - e^{-x}}{e^x + e^{-x}} = \frac{e^{2x} - 1}{e^{2x} + 1}$$

• GeLU

$$\text{GeLU}(x) = 0.5 \times x \times (1 + \tanh(\sqrt{2/\pi} \times (x + 0.044715 \times x^3)))$$

• SWiGLU

$$\text{Swish}(x) = x \cdot \text{sigmoid}(\beta x) = x / (1 + e^{-\beta x})$$

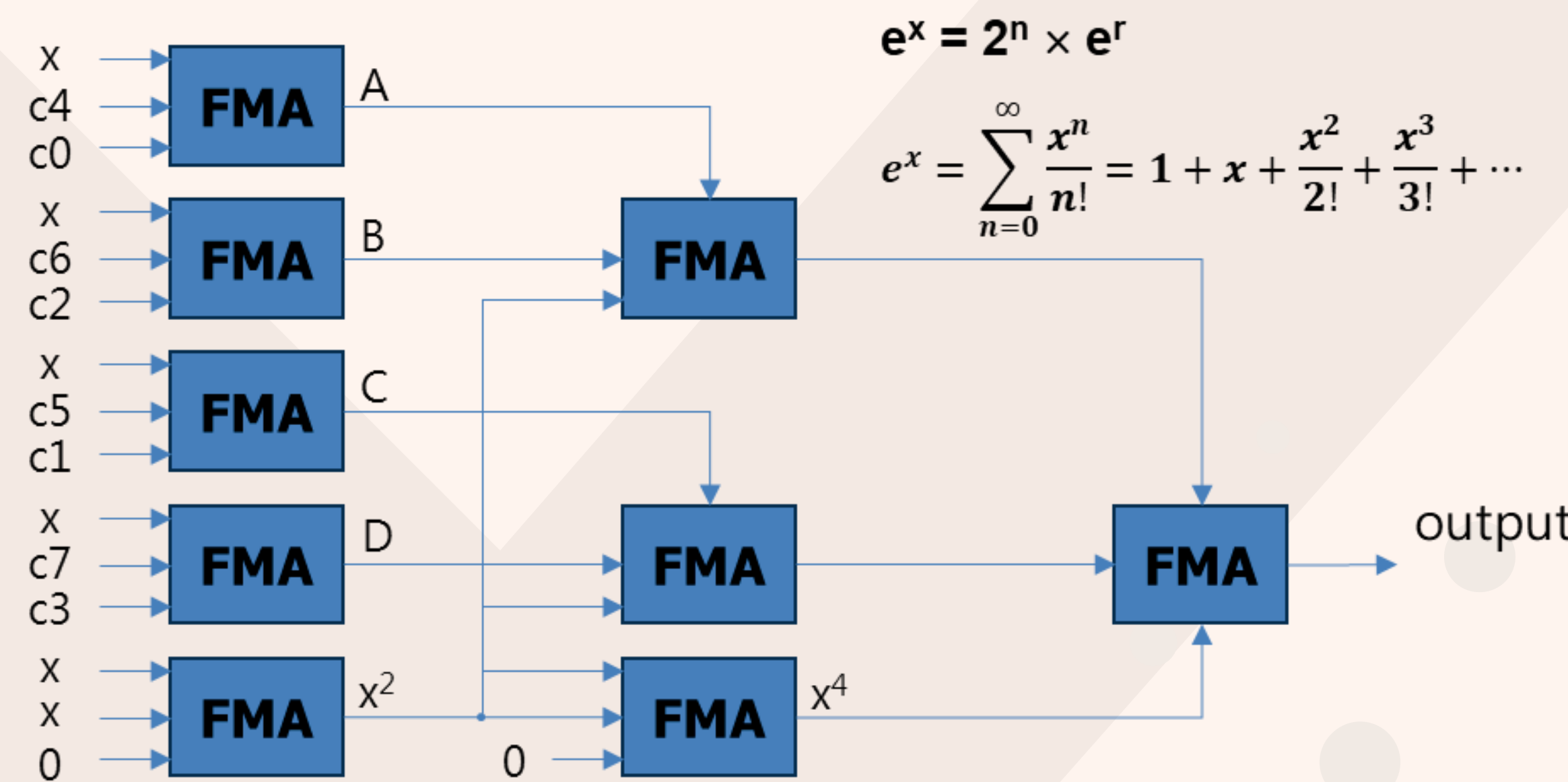
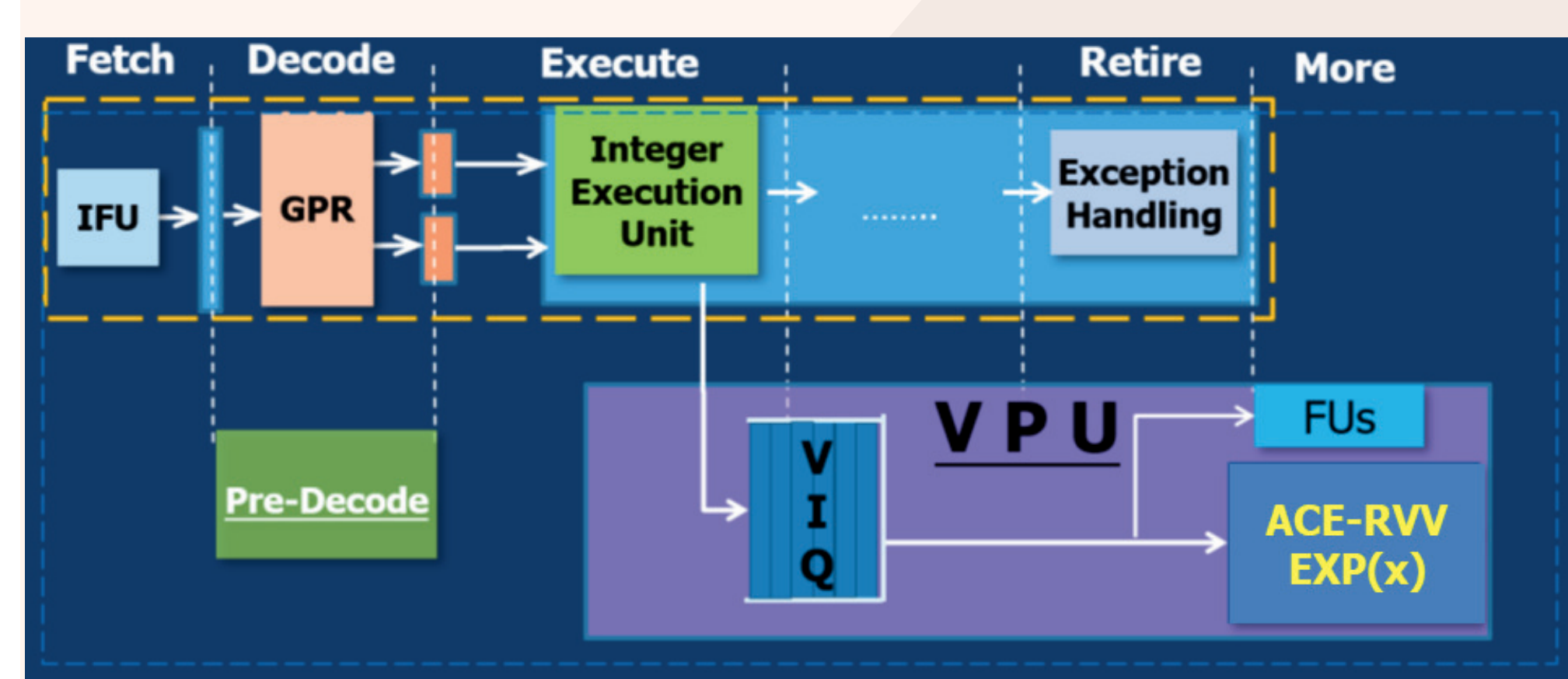
$$\text{SwiGLU}(x, W, V, b, c, \beta) = \text{Swish}_{\beta}(xW + b) \otimes (xV + c)$$

ACE-RVV Instruction Design

► Custom Exp Instruction Using ACE-RVV

► Exp Instruction by 7-Order Taylor Expansion

ACE-RVV: Reuse processor memory system



► Optimized Exp Function

Target	Throughput (Ratio)	Precision (SNR dB)
RVV + Exp	7.31	>117.08
RVV	1.00	>117.08

► Optimized Softmax Function

Target	Throughput (Ratio)
RVV + Exp	4.66
RVV	1.00

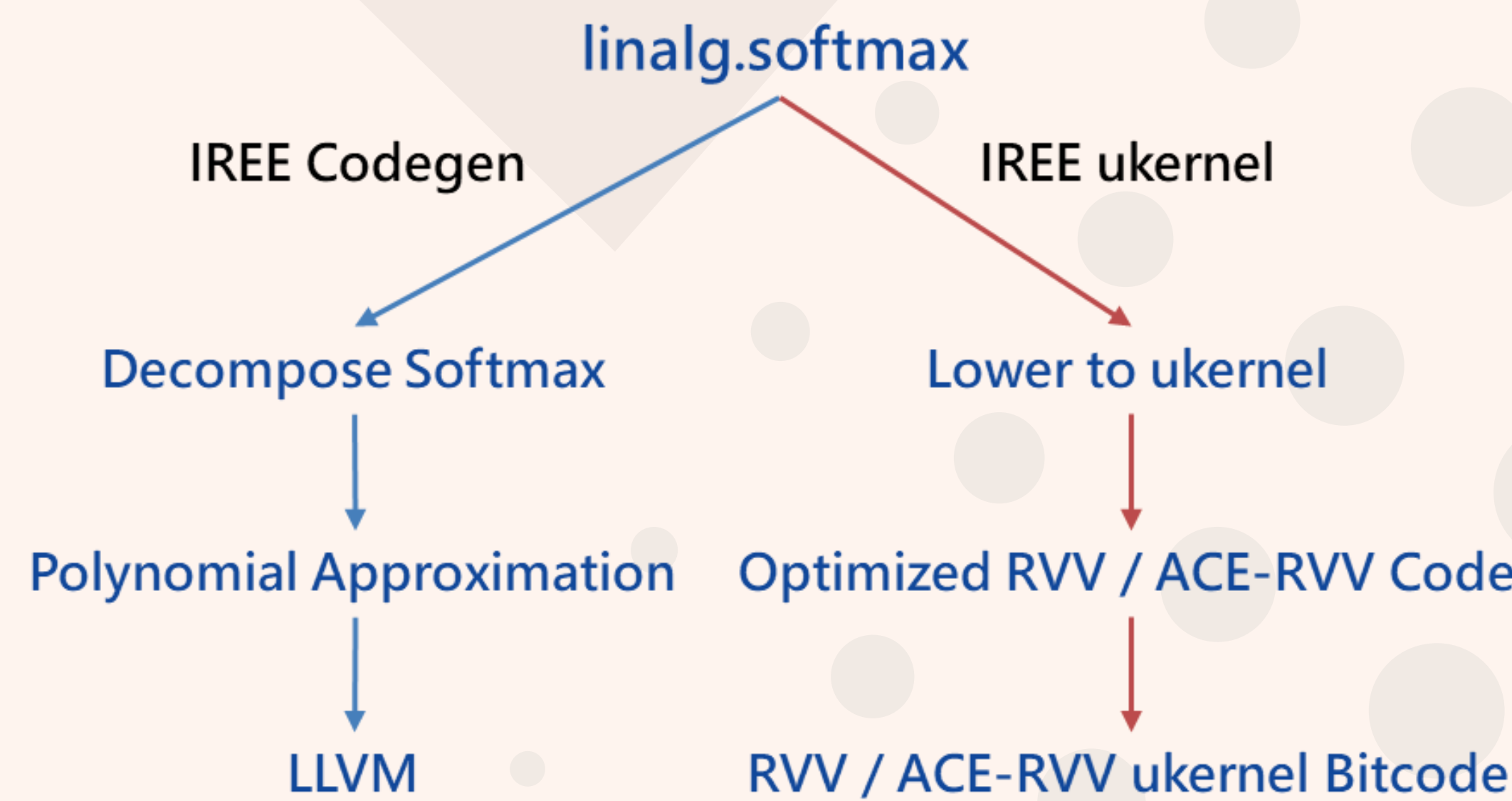
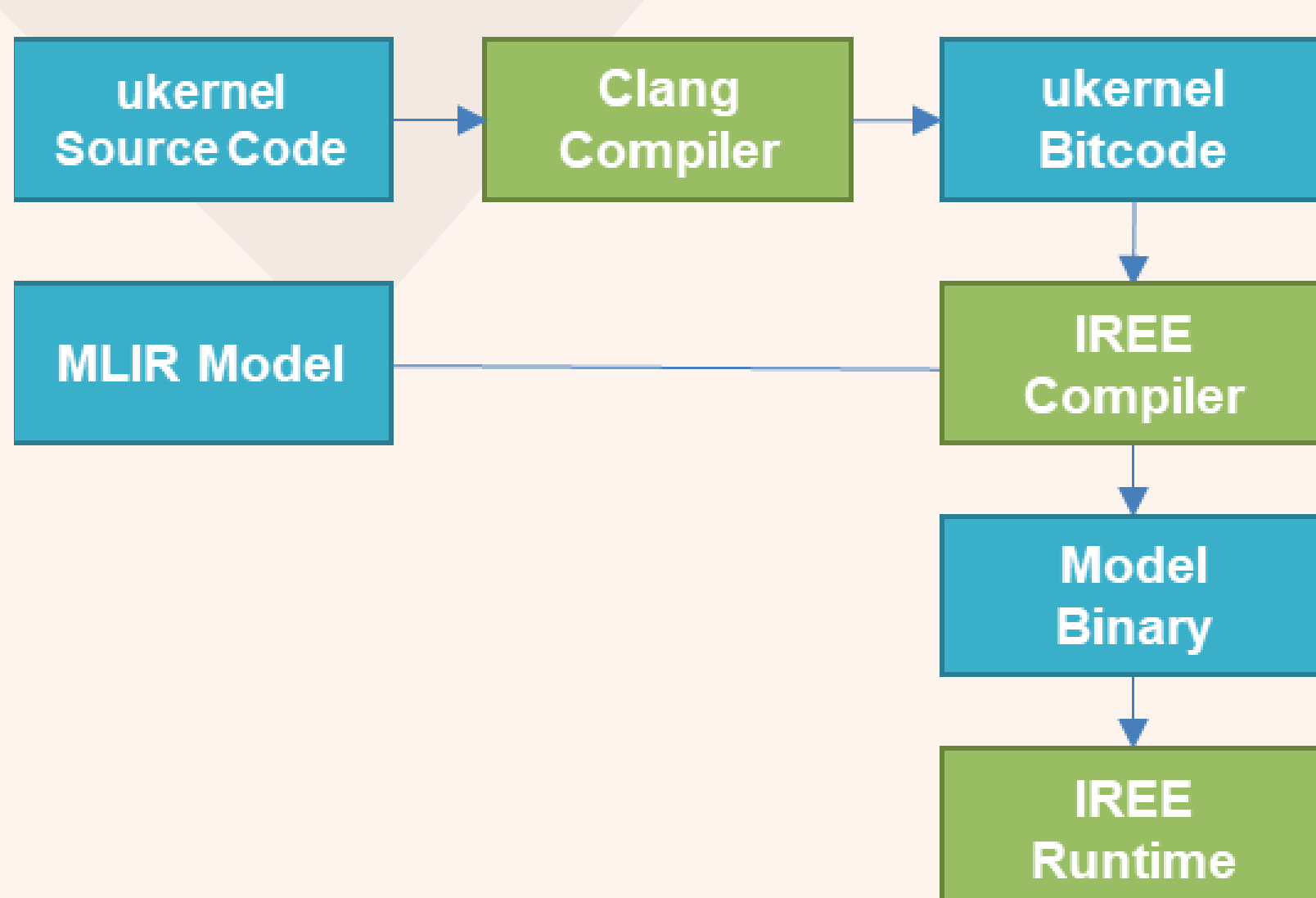
AX45MPV VLEN/DLEN = 512/512
• L1 data cache: 32KB, instruction cache: 32KB
• L2 cache: 512KB

IREE Softmax ukernel Integration

► IREE ukernel Flow

► IREE with Andes Softmax ukernel

► IREE Decompose Softmax

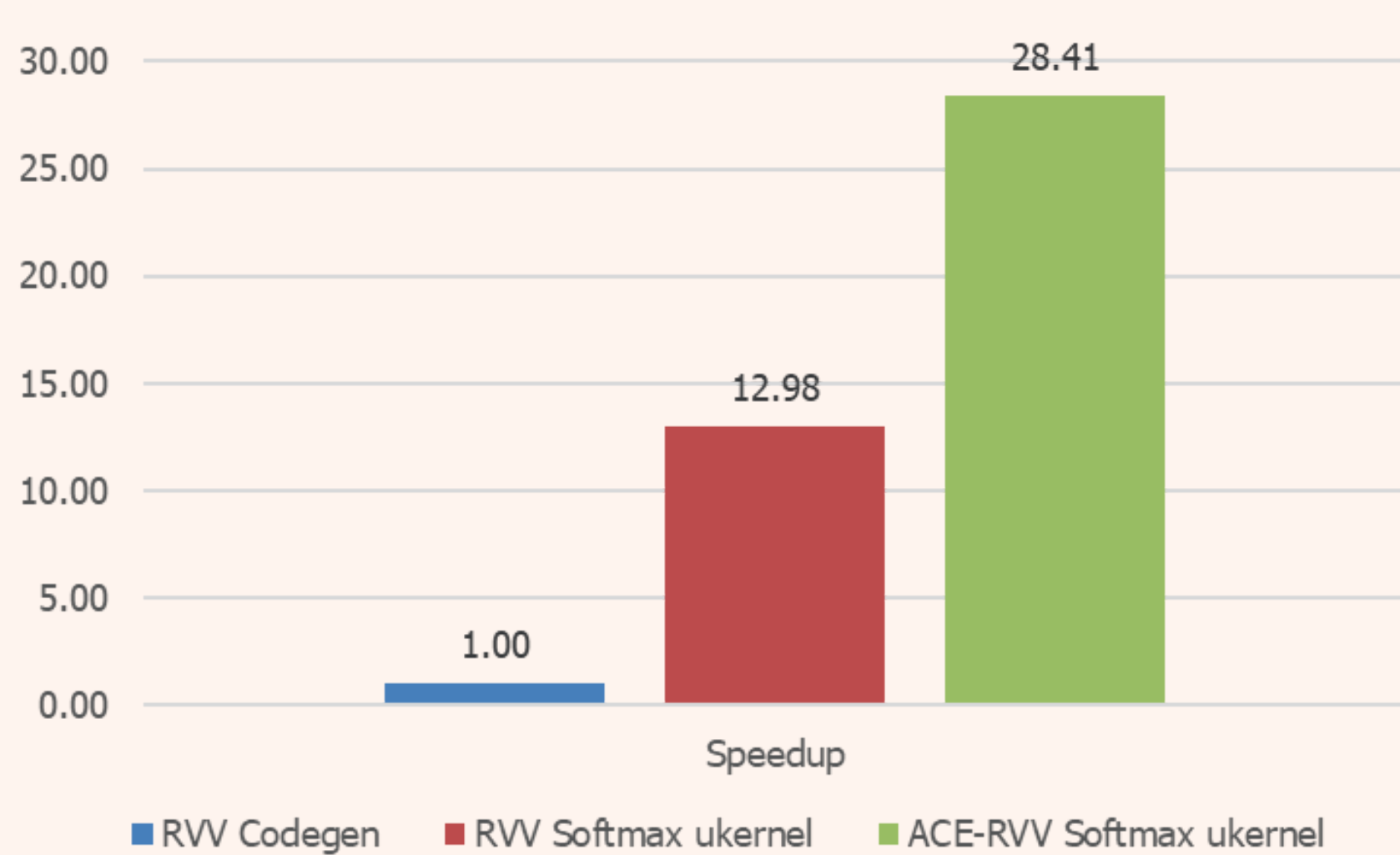


```

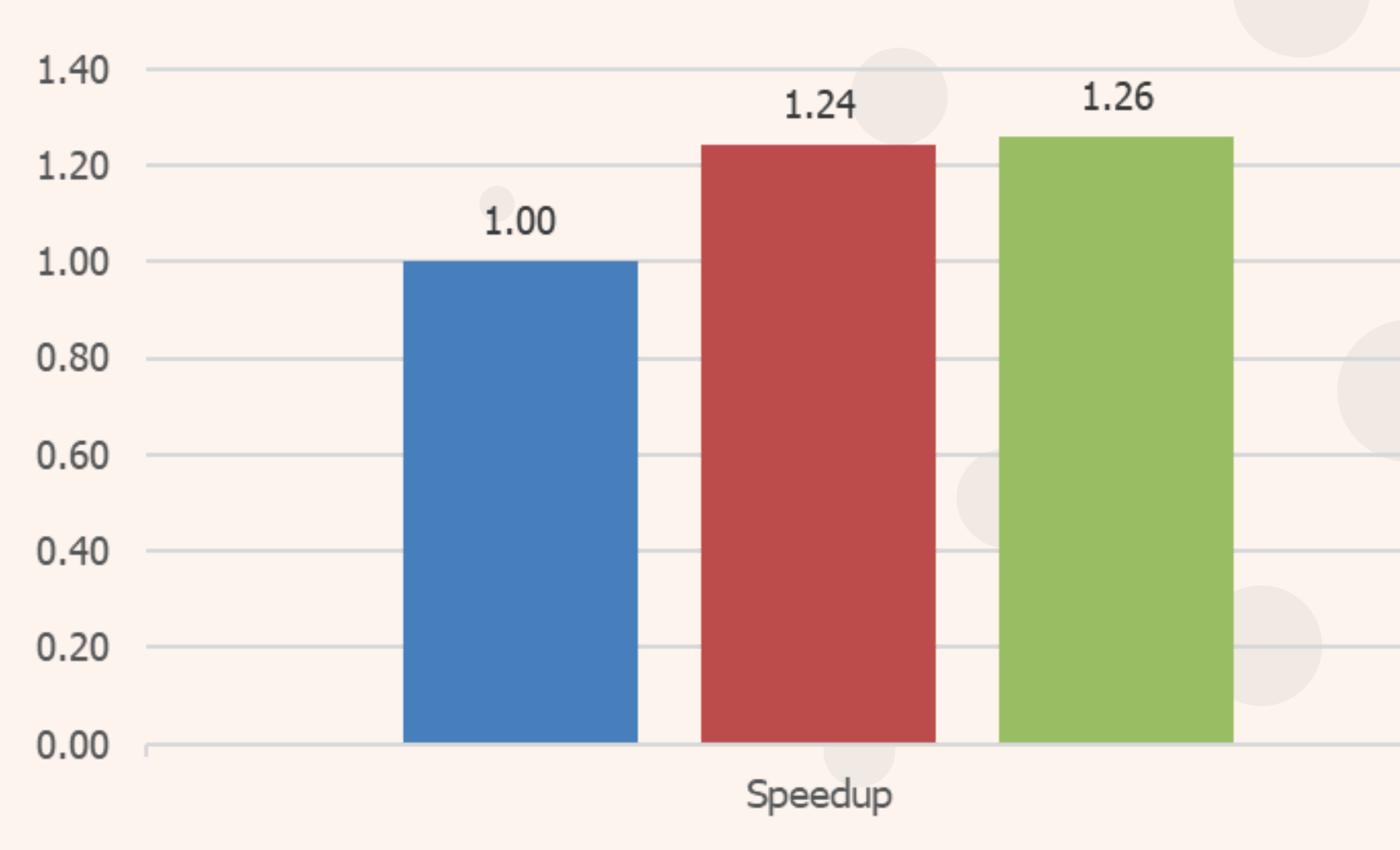
%6 = linalg.generic (indexing_maps = [affine_map<(d0, d1, d2, d3) -> (d0, d1, d2, d3)>, affine_map<(d0, d1, d2, d3) -> (d0, d1, d2)>], iterator_types = ["parallel", "parallel", "parallel", "reduction"]) ins(%2 : tensor<1x4x384x384xf32>) outs(%5 : tensor<1x4x384xf32>) {
  %bb0(@in: %2, %out: %5):
    %10 = arith.maxumf @in, %out : %32
    linalg.yield %10 : %32
  } -> tensor<1x4x384xf32>
%out_0 = arith.constant 0.000000e+00 : %32
%7 = linalg.fill ins(%out_0 : %32) outs(%4 : tensor<1x4x384xf32>) -> tensor<1x4x384xf32>
%8 = linalg.generic (indexing_maps = [affine_map<(d0, d1, d2, d3) -> (d0, d1, d2, d3)>, affine_map<(d0, d1, d2, d3) -> (d0, d1, d2)>, affine_map<(d0, d1, d2, d3) -> (d0, d1, d2)>], iterator_types = ["parallel", "parallel", "parallel", "parallel"]) ins(%2, %6, %8 : tensor<1x4x384x384xf32>, tensor<1x4x384xf32>, tensor<1x4x384x384xf32>) outs(%3, %3 : tensor<1x4x384x384xf32>):
  %bb0(@in: %2, %in_1: %3, %in_2: %3, %out: %3, %out_3: %3):
    %10 = arith.subf @in, %in_1 : %32
    %11 = math.exp %10 : %32
    %12 = arith.addf %11, %in_2 : %32
    linalg.yield %12 : %32
  } -> tensor<1x4x384xf32>
%9:2 = linalg.generic (indexing_maps = [affine_map<(d0, d1, d2, d3) -> (d0, d1, d2, d3)>, affine_map<(d0, d1, d2, d3) -> (d0, d1, d2)>, affine_map<(d0, d1, d2, d3) -> (d0, d1, d2)>], iterator_types = ["parallel", "parallel", "parallel", "parallel"]) ins(%2, %6, %8 : tensor<1x4x384x384xf32>, tensor<1x4x384xf32>, tensor<1x4x384x384xf32>) outs(%3, %3 : tensor<1x4x384x384xf32>):
  %bb0(@in: %2, %in_1: %3, %in_2: %3, %out: %3, %out_3: %3):
    %10 = arith.subf @in, %in_1 : %32
    %11 = math.exp %10 : %32
    %12 = arith.divf %11, %in_2 : %32
    linalg.yield %11, %12 : %32, %32
  } -> (tensor<1x4x384x384xf32>, tensor<1x4x384x384xf32>)
  
```

Experimental Results

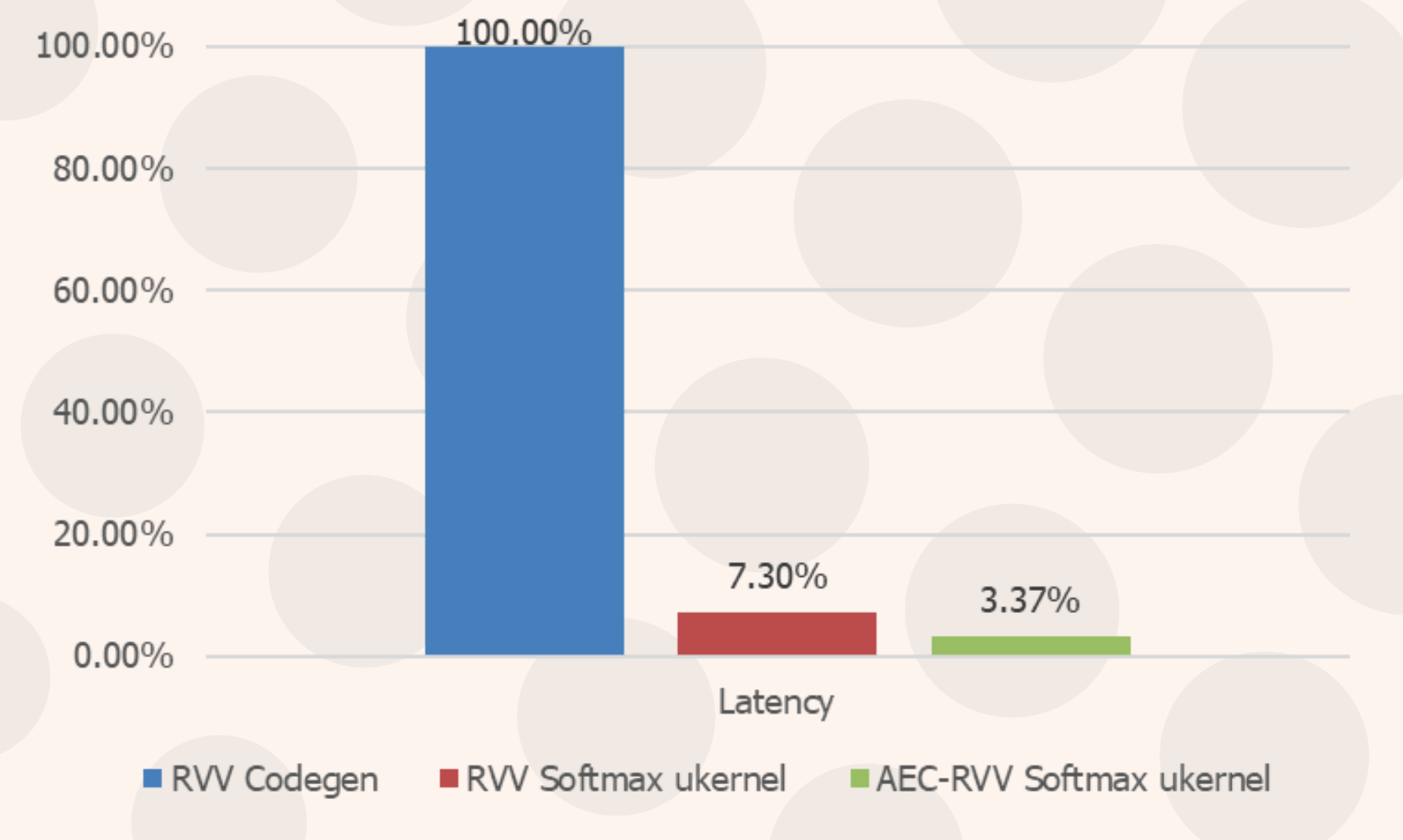
► Softmax Unit Test Speedup



► MobileBERT FP32 Speedup



► MobileBERT Softmax Op Latency



• Using IREE with Andes ACE-RVV softmax ukernel, MobileBERT has achieved a 1.26x speedup, and softmax latency has been reduced to 3.37%.